

# An Alternate Arithmetic for Word Vector Analogies

Silviu Pitis  
silviu.pitis@gmail.com

June 7, 2016

## Abstract

In recent years, word vector arithmetic of the type  $\text{king} - \text{man} + \text{woman} \approx \text{queen}$  has been applied to solve an analogy task in order to evaluate word embeddings produced by various algorithms. This paper investigates the idea of treating the relationships between word vectors as rotations of the embedding space instead of as vector differences, and shows that such treatment can produce better average approximations of target words.

## 1 Introduction

In recent years, word vector arithmetic of the type  $\text{king} - \text{man} + \text{woman} \approx \text{queen}$  has been applied to an analogy task in order to evaluate word embeddings produced by various algorithms (Mikolov et al., 2013b; Mnih and Kavukcuoglu, 2013; Pennington et al., 2014). Embeddings that do well on this task are said to capture *linguistic regularities* (Mikolov et. al) or *relational similarities* (Turney, 2006).

Such arithmetic is mysterious in the following sense: before taking the vector sum  $\text{king} - \text{man} + \text{woman}$ , the literature normalizes each vector so that it lies on the unit hypersphere (i.e.,  $\|x\| = 1$ ), but generally, the vector sum  $\text{king} - \text{man} + \text{woman}$  does not itself lie on the unit hypersphere. Although the similarity between vectors is measured by cosine distance, which ignores lengths, it seems odd that the best way to combine  $\text{king}$ ,  $\text{queen}$  and  $\text{woman}$  is to produce a vector not of length 1.

To illustrate the two dimensional case, consider *Figure 1*, where we see that the vector sum  $\text{king} - \text{man} + \text{woman}$  ends up inside the unit circle. We would need to normalize the sum to obtain an approximation for the  $\text{queen}$  unit vector, but this can alter relationships (note how  $\text{queen} - \text{woman}$  and  $\text{king} - \text{man}$  do not point in the same direction).

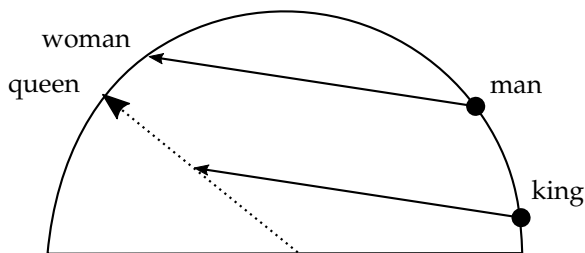


Figure 1:  $\|\text{king} - \text{man} + \text{woman}\| < 1$ .

By treating the relationship between  $\text{man}$  and  $\text{woman}$  as a rotation of the embedding space instead of a difference, I show that we can apply that same rotation to the vector for  $\text{king}$  to achieve, on average, a better vector approximation for  $\text{queen}$ .

Using pretrained word embeddings, generated by word2vec, I show that applying a rotation to the analogy task produces approximations that are closer (as determined by cosine similarity) to the target word in more than 75% of the analogies in the GOOGLE and MSR datasets. In 4 of the 14 subcategories of the GOOGLE dataset, using a rotation beats arithmetic in more than 90% of cases.

Although the *approximations* produced by applying a rotation are better on average, the *predictions* generated by selecting the closest words from the vocabulary to such approximation are mixed. When considering the single closest word, simple arithmetic beats rotations in most cases. When considering the ten closest words, however, the rotation approach is superior.

## 2 Problem Description

Let  $a, b, c, d$  be the word embedding vectors (all normalized to unit norm)<sup>1</sup> in the analogy,  $a$  is to  $b$  as  $c$  is to  $d$ . The analogy task asks us to pick the best  $d$  given  $a, b$  and  $c$ . For example,  $a, b, c, d$  might be the word embeddings for man, woman, king, queen, respectively.

Mikolov et al. (2013a) show that embeddings trained by word2vec achieve surprisingly accurate results on this task through the simple arithmetic and cosine similarity.

That is, they choose  $d$  as:

$$d = \arg \max_{d^* \in \mathcal{V}} (\text{sim}(d^*, c - a + b))$$

where

$$\text{sim}(u, v) = \frac{u \cdot v}{\|u\| \|v\|}.$$

This is odd because  $c - a + b$  is serving as a prediction for  $d$ , but  $\|d\| = 1$  and, in general,  $\|c - a + b\| \neq 1$ . We could normalize the prediction, but as can be seen in *Figure 1*, normalizing the prediction alters the direction of linear differences such as queen - woman. *Figure 2* shows the 3-dimensional case.

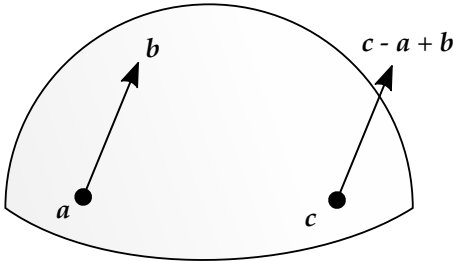


Figure 2: Simple arithmetic in 3 dimensions.

We can overcome this problem by treating the relationship between word vectors as a norm-preserving rotation of the  $d$ -dimensional embedding space,  $R : \mathbb{R}^d \rightarrow \mathbb{R}^d$  rather than as a linear difference, because  $\|R\mathbf{x}\| = \|\mathbf{x}\| \forall \mathbf{x} \in \mathbb{R}^d$ .

By treating the relationships between word vectors in this way, I show that we can, on average, improve the approximation for  $d$ , as measured by cosine similarity.

## 3 Models

In this section, I define the 3 candidate rotations to be applied to  $c$  in order to approximate  $d$ .

<sup>1</sup>The literature generally normalizes the vectors before adding them. See Mikolov et al. (2013b) and Mnih and Kavukcuoglu (2013). The author's own experiments confirmed that using unit vectors produces better results on the analogy task described.

### 3.1 Rotation from $a$ to $b$

The analogy,  $a : b :: c : d$ , says that the relationship between  $a$  and  $b$  is the same as the relationship between  $c$  and  $d$ . Thus, if  $b$  is a function,  $f$ , of  $a$ ,  $b = f(a)$ , then by hypothesis,  $d = f(c)$ .

Define  $R_{a \rightarrow b}$  as the rotation of the embedding space using the plane containing the origin,  $a$  and  $b$  as the plane of rotation that maps  $a$  onto  $b$ . It will be a rotation by  $\theta$ , where  $\cos \theta = \text{sim}(a, b) = a \cdot b$ .

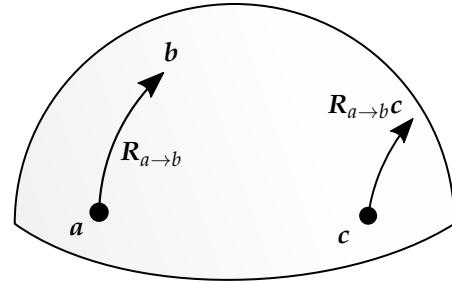


Figure 3:  $R_{a \rightarrow b}$  in 3 dimensions.

Then  $b = R_{a \rightarrow b} a$ , and we can approximate  $d$  as  $R_{a \rightarrow b} c$  and predict  $d$  as:

$$d = \arg \max_{d^* \in \mathcal{V}} (\text{sim}(d^*, R_{a \rightarrow b} c)).$$

This is illustrated in *Figure 3*.

### 3.2 Midpoint rotation

If  $a : b :: c : d$  then  $c : d :: a : b$ . But  $b = R_{a \rightarrow b} a$  whereas  $d$  is only approximated by  $R_{a \rightarrow b} c$ .

Instead, consider the rotation of the embedding space,  $R_{mid}$ , that maps the normalized midpoint of the segment  $(a, c)$  onto the normalized midpoint of the segment  $(b, d)$ . For such  $R_{mid}$ , we expect that:

$$\text{sim}(b, R_{mid} a) = \text{sim}(d, R_{mid} c).$$

This is illustrated in *Figure 4*.

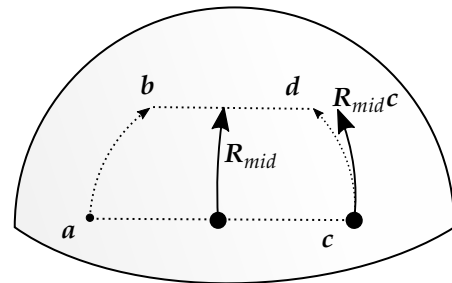


Figure 4:  $R_{mid}$  in 3 dimensions.

For the analogy task we don't know  $d$  and cannot precisely calculate  $R_{mid}$ . However, we can approximate it by substituting an approximation for  $d$ , such as  $R_{a \rightarrow b}c$ .

Therefore, we choose:

$$R_{mid} = R_{p \rightarrow q}, \text{ where}$$

$$p = \frac{a + c}{2}, \text{ and } q = \frac{b + R_{a \rightarrow b}c}{2}$$

as the rotation of the embedding space using the plane containing the origin,  $p$  and  $q$  as the plane of rotation that maps  $\frac{p}{\|p\|}$  onto  $\frac{q}{\|q\|}$ .

### 3.3 Three point rotation

A third candidate asserts that if  $a : b :: c : d$  then  $a : c :: b : d$ . The latter analogy is often nonsense (or rather arbitrary), so we wouldn't expect this to work that well, except in special cases. In any case, if this were true, then the following transformation should produce  $d$  when applied to  $c$ :

$$R_{3point} = R_{c \rightarrow a}^{-1} R_{a \rightarrow b} R_{c \rightarrow a}$$

This is illustrated in Figure 5.

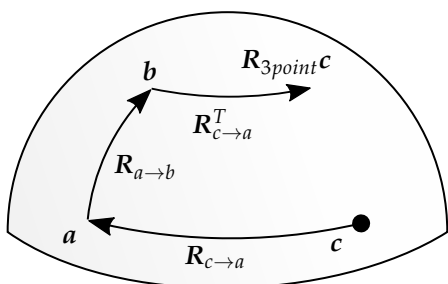


Figure 5:  $R_{3point}$  in 3 dimensions.

## 4 Experiments

### 4.1 Experiment 1: Analogy Closeness

For  $R_i \in \{R_{a \rightarrow b}, R_{mid}, R_{3point}\}$ , the first experiment compares:

$$\text{sim}(d, c - a + b)$$

to

$$\text{sim}(d, R_i c)$$

<sup>2</sup>code.google.com/archive/p/word2vec/

<sup>3</sup>word2vec.googlecode.com/svn/trunk/questions-words.txt

<sup>4</sup>research.microsoft.com/en-us/projects/rnn

<sup>5</sup>The excluded analogies were possessives like daughter : daughter's :: school : school's, since daughter's and school's were not in the vocabulary. One could drop the apostrophe and include these as plurals, but plural relationships are already included in the dataset and certain plurals like lifes would not make sense.

where  $d$  is the word vector for the true value. This experiment compares the approximation to the true value; it does not produce a prediction by considering all possible values in the vocabulary to pick the closest value, which task is performed in Section 4.3.

The word vectors used were pre-trained word2vec embeddings made available by Google<sup>2</sup>, which use a vocabulary size of 3 million and 300 dimensional embeddings.

The analogies used were taken from the GOOGLE dataset of 19544 analogies<sup>3</sup> and the MSR dataset of 8000 analogies.<sup>4</sup> Of the 8000 analogies in the full MSR dataset, 978 were skipped due to words not being in the vocabulary, leaving 7022 analogies.<sup>5</sup>

The results show that both  $R_{a \rightarrow b}$  and  $R_{mid}$  produce better approximations than simple arithmetic in a majority of cases for both the GOOGLE and MSR analogies. See Table 1 for a summary of the results, where "Avg. Diff." is defined as the average of

$$\text{sim}(d, R_i c) - \text{sim}(d, c - a + b)$$

over all analogies, and "Win %" is defined as the percentage of analogies for which

$$\text{sim}(d, R_i c) > \text{sim}(d, c - a + b).$$

It is interesting that the average cosine distance between the approximation and the true result is less for approximations produced by  $R_{a \rightarrow b}$  than for those produced by  $R_{mid}$ , but  $R_{mid}$  outperforms simple arithmetic in more cases than does  $R_{a \rightarrow b}$ . A direct comparison of  $R_{a \rightarrow b}$  to  $R_{mid}$ , using  $R_{a \rightarrow b}$  as the baseline, in Table 2, shows that  $R_{a \rightarrow b}$  is a better choice for producing approximations despite  $R_{mid}$  outperforming simple arithmetic in more cases.

### 4.2 Experiment 2: Closeness by Category

The second experiment compares the winners,  $R_{a \rightarrow b}$  and  $R_{mid}$ , of the first experiment to the baseline simple arithmetic approach across each of the 14 types of analogies in the GOOGLE dataset. The categories and results are shown below in Table 3.

In 4 of the 14 subcategories, using the rotation  $R_{mid}$  beats arithmetic in more than 90% of cases.

That there is such variance between different categories is consistent with the past work of Levy and Goldberg (2014), who show that the accuracy of simple arithmetic on the analogy task ranges from 14.55% (for the `currency` category) to 90.51% (for the `capital-common-countries` category).

### 4.3 Experiment 3: Evaluating Analogies

This experiment evaluates the effectiveness of simple arithmetic,  $R_{a \rightarrow b}$  and  $R_{mid}$  for solving analogies by finding the nearest neighbors to the approximations they produce. The experiment asks whether  $d$  is closest word in the vocabulary, as measured by cosine similarity, to the approximation produced by the model (e.g.,  $R_{a \rightarrow b}c$ ). Separately, I considered whether  $d$  was among the *ten* closest words.

Because the closest word to the approximation is often either  $b$  or  $c$ , and consistent with past literature, I removed  $b$  and  $c$  from the vocabulary when choosing the closest word. See Mnih and Kavukcuoglu (2013). I did not similarly remove  $b$  and  $c$  when looking at the ten closest words, as they are not problematic in this case.

The results, in *Table 4*, show that  $R_{a \rightarrow b}$  performs worse than simple arithmetic when considering only the closest word, but does better when considering the ten closest words.  $R_{mid}$  was competitive with simple arithmetic when looking only at the closest word, and outperformed simple arithmetic when looking at the ten closest words.

Although the accuracies of all methods improved when given ten guesses, it is interesting that the improvement of  $R_{a \rightarrow b}$  is the sharpest; allowing extra guesses lifts  $R_{a \rightarrow b}$  from last place on both datasets to being the winner of the MSR dataset and the most categories in the GOOGLE dataset.

## 5 A Note on SemEval-2012 Task 2

SemEval-2012 Task 2,<sup>6</sup> appearing in Jurgens et al. (2012), has been used to measure semantic regularities in word vectors (Mikolov et al., 2013b; Levy and Goldberg, 2014). The SEMEVAL dataset covers 79 categories of semantic relationships. For each category, 3 or 4 prototypical word pairs exemplify the relationship, and the task involves ranking approximately 40 target word pairs in that category according to how well they reflect the relationship. For example, the `CLASS-INCLUSION: Taxonomic category` is exemplified by the word pairs, `flower:tulip`,

`emotion:rage` and `poem:sonnet`, and the task is to rank 41 target word pairs including `hair:brown` and `pet:dog` by the degree to which they belong in the same category.

If  $a : b$  is a prototypical word pair and  $c : d$  is a target word pair, past authors have evaluated the relational similarity in two ways:

1. as an analogy, by evaluating  $sim(d, c - a + b)$ , or
2. as a relationship, by evaluating  $sim(d - c, b - a)$ .

Mikolov et al. (2013b) and Levy and Goldberg (2014) found that the latter option produces better results.

Thinking in terms of rotations changes the former option. That is, instead of evaluating  $sim(d, c - a + b)$ , we can evaluate  $sim(d, Rc)$ . This produces the following results:

	Arithmetic	$R_{a \rightarrow b}$	$R_{mid}$
Accuracy	40.7%	37.7%	<b>41.1%</b>

The second option, comparing the relationships directly, achieves an accuracy of **44.8%**, outperforming all analogical approaches. Thinking in terms of rotations becomes a bit trickier for this option. We can calculate each relationship as  $R_{a \rightarrow b}$  and  $R_{c \rightarrow d}$ , but there is no best way to determine their similarity.

One approach is to use a generalized concept of an angle that allows an angle to be defined between two arbitrary subspaces and then rank target word pairs by the angle between the planes of rotation of  $R_{a \rightarrow b}$  and  $R_{c \rightarrow d}$ . This throws out the direction of the rotation, but we can correct this by multiplying by the sign of  $sim(d - c, b - a)$ . Applying the definition of an angle between flats given by Jordan, the algorithm for which is described in Knyazev and Argentati (2002) and implemented in the `krypy`<sup>7</sup> python package, this approach achieves a relatively respectable accuracy of **43.5%**.

The easier, and perhaps more natural, approach to comparing the rotations is precisely the approach of Mikolov et al. (2013b) described above. The difference vectors  $b - a$  and  $d - c$  lie in the planes of rotation of  $R_{a \rightarrow b}$  and  $R_{c \rightarrow d}$ , respectively. In addition to the plane of rotation, the difference vectors also contain information about the *direction* and *magnitude* of the rotation.

<sup>6</sup><https://sites.google.com/site/semEval2012task2/>

<sup>7</sup><https://github.com/andrenarchy/krypy>

Category	No. of Analogies	$R_{a \rightarrow b}$		$R_{mid}$	
		Avg. Diff.	Win %	Avg. Diff.	Win %
capital-common-countries	506	<b>0.023</b>	61.7%	0.022	<b>83.8%</b>
capital-world	4524	<b>0.035</b>	78.0%	0.022	<b>85.2%</b>
city-in-state	2467	<b>0.022</b>	<b>67.2%</b>	-0.003	50.6%
currency	866	-0.000	<b>50.2%</b>	-0.002	47.6%
family	506	<b>0.028</b>	<b>80.2%</b>	0.011	72.5%
gram1-adjective-to-adverb	992	<b>0.067</b>	<b>85.0%</b>	0.019	76.1%
gram2-opposite	812	<b>0.050</b>	<b>75.6%</b>	0.005	60.6%
gram3-comparative	1332	-0.007	40.0%	<b>0.006</b>	<b>62.5%</b>
gram4-superlative	1122	-0.004	45.7%	<b>0.006</b>	<b>59.1%</b>
gram5-present-participle	1056	<b>0.067</b>	92.4%	0.047	<b>98.6%</b>
gram6-nationality-adjective	1599	0.007	51.8%	<b>0.017</b>	<b>88.7%</b>
gram7-past-tense	1560	<b>0.063</b>	87.2%	0.044	<b>94.1%</b>
gram8-plural	1332	<b>0.080</b>	95.4%	0.049	<b>96.8%</b>
gram9-plural-verbs	870	<b>0.051</b>	75.9%	0.039	<b>90.1%</b>

Table 3: Comparison of  $R_{a \rightarrow b}$  and  $R_{mid}$  to simple arithmetic by analogy category in the **GOOGLE** dataset

Category	Closest word			Closest ten words		
	Arithmetic	$R_{a \rightarrow b}$	$R_{mid}$	Arithmetic	$R_{a \rightarrow b}$	$R_{mid}$
GOOGLE dataset	<b>54.1%</b>	41.0%	52.5%	60.7%	61.2%	<b>62.6%</b>
MSR dataset	39.1%	30.3%	<b>40.4%</b>	42.5%	<b>48.7%</b>	47.9%
GOOGLE dataset by category						
capital-common-countries	<b>61.5%</b>	49.8%	60.9%	67.0%	62.1%	<b>68.8%</b>
capital-world	<b>62.4%</b>	43.4%	61.5%	68.5%	68.0%	<b>71.3%</b>
city-in-state	<b>41.1%</b>	27.1%	35.8%	<b>46.3%</b>	<b>46.3%</b>	43.5%
currency	<b>26.8%</b>	4.5%	21.6%	<b>38.0%</b>	17.3%	33.9%
family	<b>75.1%</b>	62.6%	72.9%	<b>86.0%</b>	84.6%	85.4%
gram1-adjective-to-adverb	<b>18.0%</b>	13.3%	17.6%	27.2%	<b>38.5%</b>	27.8%
gram2-opposite	<b>31.2%</b>	20.1%	26.8%	<b>37.2%</b>	31.2%	33.7%
gram3-comparative	56.2%	38.8%	<b>57.7%</b>	56.7%	<b>63.5%</b>	60.1%
gram4-superlative	<b>36.3%</b>	15.0%	31.4%	38.2%	35.1%	<b>42.6%</b>
gram5-present-participle	62.4%	62.6%	<b>64.2%</b>	75.4%	<b>83.8%</b>	82.1%
gram6-nationality-adjective	86.1%	82.9%	<b>86.2%</b>	92.3%	92.5%	<b>92.9%</b>
gram7-past-tense	<b>49.0%</b>	37.9%	47.6%	60.6%	<b>65.9%</b>	65.4%
gram8-plural	74.5%	78.0%	<b>79.4%</b>	77.0%	<b>85.5%</b>	83.9%
gram9-plural-verbs	<b>49.7%</b>	20.0%	41.8%	59.5%	50.3%	<b>61.5%</b>

Table 4: Accuracy on analogy task at precision levels of 1 and 10.

Therefore, the approach taken by Mikolov et al. (2013b) and Levy and Goldberg (2014) to SemEval-2012 Task 2 is consistent with the idea of thinking about relationships between word vectors as rotations of the embedding space.

## 6 Conclusion

I've introduced a new way to evaluate and think about the relationships between word vectors: as a rotation of the embedding space instead of as a vector difference. So long as we are using cosine similarity to evaluate the similarity between normalized word vectors, this method makes better geometric sense. When used for the word analogy task, it produces a better average approximation of the target word vector.

These results are neither groundbreaking nor immediately useful from a practical perspective. They do, however, offer a new way to think about and visualize the word embedding space, which may be helpful for future research. Avenues for follow-up may include exploring whether applying rotations to the original unnormalized word embedding can produce better unnormalized approximations, and the use of vector approximations to expand a trained vocabulary with new or foreign words.

## References

- Jurgens, D., Mohammad, S., Turney, P., and Holyoak, K. (2012). Semeval-2012 task 2: Measuring degrees of relational similarity. *First Joint Conference on Lexical and Computational Semantics (\*SEM)*, pages 356–364.
- Knyazev, A. V. and Argentati, M. E. (2002). *Principal Angles Between Subspaces in An A-Based Scalar Product: Algorithms and Perturbation Estimates*, volume 23.
- Levy, O. and Goldberg, Y. (2014). Linguistic regularities in sparse and explicit word representations. In *Proceedings of the Eighteenth Conference on Computational Natural Language Learning*, pages 171–180. Association for Computational Linguistics.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a). Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781.
- Mikolov, T., tau Yih, W., and Zweig, G. (2013b). Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for*

*Computational Linguistics: Human Language Technologies (NAACL-HLT-2013)*. Association for Computational Linguistics.

- Mnih, A. and Kavukcuoglu, K. (2013). Learning word embeddings efficiently with noise-contrastive estimation. In Burges, C. J. C., Bottou, L., Welling, M., Ghahramani, Z., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems 26*, pages 2265–2273. Curran Associates, Inc.
- Pennington, J., Socher, R., and Manning, C. D. (2014). GloVe: Global Vectors for Word Representation. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*, pages 1532–1543.
- Turney, P. D. (2006). Similarity of Semantic Relations. *Computational Linguistics*, (February):1–39.